

Distributed Storage Systems based on Equidistant Subspace Codes

Netanel Raviv

Tuvi Etzion

June 25, 2014

Abstract

Distributed storage systems based on equidistant constant dimension codes are presented. These equidistant codes are based on the Plücker embedding, which is essential in the repair and the reconstruction algorithms. These systems possess several useful properties such as high failure resilience, minimum bandwidth, low storage, simple algebraic repair and reconstruction algorithms, good locality, and compatibility with small fields.

1 Introduction

Let q be a prime power and let \mathbb{F}_q be the field with q elements. In a distributed storage system (DSS) a file $x \in \mathbb{F}_q^B$ is stored in n storage *nodes*, α information symbols in each. The DSS is required to be resilient to node failures; i.e., it should be possible to retrieve the data from a lost node by contacting d other active nodes and downloading β information symbols from each one of them, an operation which is called *repair*. In addition, a *data collector* (DC) should be able to rebuild the stored file x by contacting k active nodes, an operation which is called *reconstruction*. If the file is coded with an ordinary error correcting code C prior to being stored in the system (usually by an MDS code [1, 2, 3, 4, 5, 6, 7, 8]), then C is called the *outer code*, and the DSS code is called the *inner code*.

A repair process that results in a new node which contains the exact same information as in the failed node is called an *exact repair* [1, 9]. A repair process which is not an exact repair is called a *functional repair*. Such a repair must maintain the system's ability of repair and reconstruction. The amount of data which is required for a repair is $d\beta$, and it is called the *repair bandwidth* of the code. Codes which minimize the repair bandwidth, i.e., $d\beta = \alpha$, are called Minimum Bandwidth Regenerating (MBR) Codes [10]. Codes which minimize α , and thus have $\alpha = \frac{B}{k}$, are called Minimum Storage Regenerating (MSR) Codes [10]. A Self-Repairing Code (SRC) [11] is a code satisfying: (a) repairs are possible without having to download an amount of data equivalent to the reconstruction of the original file x ; and (b) the number of nodes required for repair depends only on how many nodes are missing and not on their identity.

In [12] a framework for a construction of a DSS code based on subspaces is given. This framework is slightly different from the classical one. In this framework every node v_i is associated with a

¹This research was supported in part by the Israeli Science Foundation (ISF), Jerusalem, Israel, under Grant 10/12. The work of Netanel Raviv is part of his Ph.D. thesis performed at the Technion. The authors are with the Department of Computer Science, Technion, Haifa 32000, Israel. e-mail: etzion,netanel@cs.technion.ac.il .

subspace U_i of a vector space U called the *message space*. The dimension of U is $B = |x|$, where $x \in \mathbb{F}_q^B$ is the file to be stored. In the “storage phase” a node v_i receives a vector $M_i \cdot x$, where M_i is a full-rank matrix whose row span is U_i . A set of nodes is called a *reconstruction set*¹ if their respective subspaces span the entire message space. The file x is reconstructible from a reconstruction set $\{v_i\}_{i \in I}, I \subseteq [n]$, where $[\ell] \triangleq \{1, \dots, \ell\}$, by solving a linear nonsingular equation system based on $\{M_i \cdot x\}_{i \in I}$ and $\{M_i\}_{i \in I}$. A set $\{v_i\}_{i \in T_j}, T_j \subseteq [n]$ of nodes is called a *repair set for a node v_j* if each subspace $U_i, i \in T_j$ contains a subspace $W_{i,j} \subseteq U_i$ such that the span of the set $\{W_{i,j} \mid i \in T_j\}$ contains U_j . The lost information $M_j \cdot x$ may be retrieved by manipulating the rows in a linear system based on $\{M_{i,j} \cdot x\}_{i \in T_j}$ and $\{M_{i,j}\}$, where $M_{i,j}$ is a matrix whose row span is $W_{i,j}$. This framework yields an algebraic repair and reconstruction algorithms. We will use the equidistant subspace codes from [13] as the subspaces in our DSS. We note that in this new framework the matrices M_i have the role of the outer code in the classic framework.

Our codes achieve the SRC property, and nearly achieve the MSR and MBR properties. Regarding the MBR property, we show that $d\beta \leq \alpha + 1$, and hence the MBR property is achieved up to an additive constant of 1. Regarding the MSR property, we show that if the nodes participating in the reconstruction algorithm receive some information from the DC, then it is possible to reconstruct x by communicating $|x| = B$ field elements, $\frac{b}{2}$ elements from each node if b is even and either $\frac{b-1}{2}$ elements or $\frac{b+1}{2}$ elements if b is odd. This property may be seen as a variant of the MSR property. Without this additional assumption it is possible to reconstruct x by downloading $2B$ elements from b nodes. The penalty of providing these advantages is not being able to repair (resp. reconstruct) from *any* set of d (resp. k) nodes, but rather some properly chosen ones. This drawback is also apparent in some existing DSS codes [2, 11].

Our code stores a file $x \in \mathbb{F}_q^B$, where $B = \binom{b}{2}$ for some $b \in \mathbb{N}$, in n nodes. The user may choose any n such that $b \leq n \leq \frac{q^b - 1}{q - 1}$ in correspondence with the expected number of simultaneous node failures. Each node stores $b - 1$ field elements. For the purpose of repair, the user may choose one of two possible algorithms. The first one requires that the *newcomer node* (newcomer, in short) will contact either $b - 1$ or b active nodes and download a single field element from each one. This algorithm will minimize the *repair bandwidth* as possible. The second algorithm requires downloading all data from as little as two nodes, depending on the code construction. In either of the algorithms it is not possible to contact *any* set of nodes, but a proper set may be easily found, and it is promised to exist as long as the number of node failures does not exceed some reasonable bound.

The presented code has several useful properties. As mentioned earlier the user may choose between a local repair (Subsection 3.3) and a minimum bandwidth repair (Subsection 3.2). In addition, it is possible to reconstruct nodes that were not previously in the system (Corollary 1); that is, once a proper set of b nodes is stored in the system by the user, the system may use repairs in order to generate additional storage nodes without any outside interference. It is also possible to repair in the presence of up to $O(\sqrt{B})$ simultaneous node failures, while imposing no restriction on the field size (Example 3). Two additional useful properties are apparent. One is the ability to efficiently reuse the system to store a file $y \neq x$, without having to initialize all nodes (Subsection 3.6). This property follows directly from the linear nature of our code. The second is the ability to simultaneously repair multiple node failures in parallel (Subsection 3.4).

A brief overview of the equidistant subspace codes from [13] will be given in Section 2. The specific properties of our code strongly depend on an assignment of different vectors as identifiers to the storage nodes. The code will first be described with respect to a general assignment in Section 3, and specific assignments, as well as their resulting properties, will be discussed in Section 4. Some proofs and further explanations in this version are omitted and will appear in the full version of

¹[12] uses the term *recovery set*. We use a different term for consistency.

this paper.

2 Preliminaries

The Grassmannian $\mathcal{G}_q(n, k)$ is the set of all k -subspaces of \mathbb{F}_q^n . The size of $\mathcal{G}_q(n, k)$ is given by the Gaussian coefficient $\begin{bmatrix} n \\ k \end{bmatrix}_q$ (see [14, Chapter 24]). A *constant dimension code* (CDC) is a subset of $\mathcal{G}_q(n, k)$ with respect to the *subspace metric* $d_S(U, V) = \dim U + \dim V - 2 \dim(U \cap V)$. A CDC is called equidistant if the distance between every two distinct codewords is some fixed constant. An equidistant CDC is also called a *t-intersecting* code since the dimension of the intersection of any two distinct codewords is some constant t . Our construction uses the 1-intersecting equidistant subspace codes from [13], whose construction and properties are hereby described.

In what follows e_i denotes the i th unit vector. For a set S of vectors, $\langle S \rangle$ denotes the linear span of S , and for a matrix M , $\langle M \rangle$ denotes its row linear span.

Definition 1. (The Plücker embedding, see [13, Section 4], [15, p. 165]) Given $M \in \mathbb{F}_q^{t \times b}$, identify the coordinates of $\mathbb{F}_q^{\binom{b}{t}}$ with all t -subsets of $[b]$, and define $\varphi(M)$ as a vector of length $\binom{b}{t}$ in which

$$(\varphi(M))_{\{i_1, \dots, i_t\}} \triangleq \det M(i_1, \dots, i_t)$$

where $M(i_1, \dots, i_t)$ is the $t \times t$ sub-matrix of M formed from columns $i_1 < \dots < i_t$. For $v, u \in \mathbb{F}_q^b$ we denote by $\varphi\binom{v}{u}$ the result of applying φ on the $2 \times b$ matrix $\binom{v}{u}$.

Definition 2. [13, Subsection 3.1] For $V \in \mathcal{G}_q(b, 1)$, $v \in V \setminus \{0\}$, and the index $r(v)$ of the leftmost nonzero entry of v , let

$$P_V \triangleq \left\langle \left\{ \varphi\binom{v}{e_i} \right\}_{i \in [b] \setminus \{r(v)\}} \right\rangle.$$

By the properties of the determinant function, any choice of a nonzero vector v from the 1-subspace V results in the same subspace, and thus P_V is well-defined. Lemma 3 which follows shows that the choice of $r(v)$ as the leftmost nonzero entry of v is arbitrary, and every other nonzero entry could equally be chosen.

Theorem 1. [13, Theorem 14] The following code

$$\mathbb{C} \triangleq \{P_V \mid V \in \mathcal{G}_q(b, 1)\},$$

$\mathbb{C} \subseteq \mathcal{G}_q\left(\binom{b}{2}, b-1\right)$ is an equidistant 1-intersecting code of size $\begin{bmatrix} b \\ 1 \end{bmatrix}_q$; that is, any distinct $P_U, P_V \in \mathbb{C}$ satisfy $\dim(P_U \cap P_V) = 1$. In addition, for every distinct $P_U, P_V \in \mathbb{C}$, $P_U \cap P_V = \langle \varphi\binom{u}{v} \rangle$, where $U = \langle u \rangle$ and $V = \langle v \rangle$.

The following lemma shows that the function φ from Definition 1 is a bilinear form when applied on two row matrices. This fact will be prominent in our constructions.

Lemma 1. [13, Lemma 4] If $v, u \in \mathbb{F}_q^b$ are nonzero vectors, and $\gamma, \delta \in \mathbb{F}_q$, then $\varphi\binom{v}{\gamma u + \delta w} = \gamma \cdot \varphi\binom{v}{u} + \delta \cdot \varphi\binom{v}{w}$ and $\varphi\binom{\gamma u + \delta w}{v} = \gamma \cdot \varphi\binom{u}{v} + \delta \cdot \varphi\binom{w}{v}$.

Lemma 2 and Lemma 3 provide a convenient way of choosing a basis to any $P_V \in \mathbb{C}$ (Theorem 1); and both may easily be obtained from [13, Lemma 3]. For completeness we include a short proof.

Lemma 2. If $v = (\gamma_1, \dots, \gamma_b) \in \mathbb{F}_q^b$ is a nonzero vector, then

$$\sum_{j \in [b]} \gamma_j \cdot \varphi \begin{pmatrix} v \\ e_j \end{pmatrix} = 0.$$

Proof. By Lemma 1 and by the properties of the determinant function we have

$$\sum_{j \in [b]} \gamma_j \cdot \varphi \begin{pmatrix} v \\ e_j \end{pmatrix} = \varphi \begin{pmatrix} v \\ \sum_{j \in [b]} \gamma_j e_j \end{pmatrix} = \varphi \begin{pmatrix} v \\ v \end{pmatrix} = 0.$$

□

Lemma 3. If $v = (\gamma_1, \dots, \gamma_b) \in \mathbb{F}_q^b$, $\langle v \rangle \triangleq V$, and $\gamma_s \neq 0$ for some $s \in [b]$, then

$$P_V = \left\langle \left\{ \varphi \begin{pmatrix} v \\ e_i \end{pmatrix} \right\}_{i \in [b] \setminus \{s\}} \right\rangle.$$

Proof. By Lemma 2,

$$\begin{aligned} \varphi \begin{pmatrix} v \\ e_{r(v)} \end{pmatrix} &\in \left\langle \left\{ \varphi \begin{pmatrix} v \\ e_i \end{pmatrix} \right\}_{i \in [b] \setminus \{r(v)\}} \right\rangle \\ \varphi \begin{pmatrix} v \\ e_s \end{pmatrix} &\in \left\langle \left\{ \varphi \begin{pmatrix} v \\ e_i \end{pmatrix} \right\}_{i \in [b] \setminus \{s\}} \right\rangle, \end{aligned}$$

and hence,

$$P_V \triangleq \left\langle \left\{ \varphi \begin{pmatrix} v \\ e_i \end{pmatrix} \right\}_{i \in [b] \setminus \{r(v)\}} \right\rangle = \left\langle \left\{ \varphi \begin{pmatrix} v \\ e_i \end{pmatrix} \right\}_{i \in [b]} \right\rangle = \left\langle \left\{ \varphi \begin{pmatrix} v \\ e_i \end{pmatrix} \right\}_{i \in [b] \setminus \{s\}} \right\rangle$$

□

The following observation will be repeatedly used throughout our algorithms.

Observation 1. Let $A, B \in \mathbb{F}_q^{b \times B}$ be two distinct row-equivalent matrices. If r_1, \dots, r_t is the series of row operations that transform M_1 to M_2 , then for any $x \in \mathbb{F}_q^B$ it is possible to compute $M_2 x$ given $M_1 x$ and r_1, \dots, r_t .

Proof. Let E_1, \dots, E_t be the invertible matrices corresponding to the row operations that transform M_1 to M_2 ; that is, $E_1 \cdot E_2 \cdot \dots \cdot E_t \cdot M_1 = M_2$. The claim follows directly from the fact that for any $x \in \mathbb{F}_q^B$, $E_1 \cdot E_2 \cdot \dots \cdot E_t \cdot M_1 x = M_2 x$. □

Remark 1. The complexity analysis of Algorithms 1 through 4 in the sequel, relies mostly on the complexity of solving a system of linear equations over a finite field. This can be done either by a school book Gaussian elimination or by employing one of many faster algorithms (see [16] and references therein). However, to simplify the discussion we analyze our algorithms by using simple Gaussian elimination.

3 The Distributed Storage System

We are now in a position to describe the construction of the DSS. The feasibility of the described repair and reconstruction algorithms will depend on a certain assignment of vectors in \mathbb{F}_q^b to identify the storage nodes. Different assignments and their resulting parameters will be discussed separately in Section 4. With respect to a certain assignment of vectors to nodes, we will say that a set of nodes are *linearly independent* if their assigned vectors are linearly independent.

3.1 Storage

Let v_1, \dots, v_n be the available storage nodes. We identify each v_i by a *normalized* vector from \mathbb{F}_q^b ; that is, a vector whose leftmost nonzero entry $r(v_i)$ is 1. Let M_{v_i} be the $(b-1) \times B$ matrix whose rows are the vectors

$$\left\{ \varphi \begin{pmatrix} v_i \\ e_j \end{pmatrix} \right\}_{e_j \in [b] \setminus r(v_i)}.$$

Following the terminology in [12, Section III.A.], each node v_i is in fact associated with a subspace. In our system, this subspace is $P_{\langle v_i \rangle} \triangleq \langle M_{v_i} \rangle$ (see Definition 2).

Let s be the source node, i.e. the node holding the file $x \in \mathbb{F}_q^B$ to be stored. For the initial storage, s sends $M_{v_i} \cdot x$ to v_i for all $i = 1, \dots, n$. It is evident that $n \cdot (b-1)$ field elements are being sent. As for time complexity, computing the product $M_{v_i} \cdot x$ requires computing the matrix M_{v_i} . If the vector v_i is given, each $\varphi \begin{pmatrix} v_i \\ e_j \end{pmatrix}$ is computable from v_i in $O(b \log b)$ time by using a proper sparse representation². Hence, the matrix M_{v_i} is computable in $O(b^2 \cdot \log b) = O(B \log B)$. Using the same sparse representation, computing the product $M_{v_i} \cdot x$ takes an additional $O(B \log B)$ time for each v_i . This stage requires $O(B \log B \cdot n)$ computation time and $O(B^{1/2} \cdot n)$ communication units.

3.2 Minimum Bandwidth Repair

In what follows we show that it is possible to repair a node failure by communicating a single field element from either $b-1$ or b nodes. For functional repair no further computations are needed while for exact repair an additional $O(B^2)$ algorithm should be applied by the newcomer.

Algorithm 1. Let $v_j = \sum_{t=1}^b \gamma_t e_t$ be the failed node and let $u_1, \dots, u_{b'}$ be any set of active nodes such that $\langle e_t \rangle_{t \in [b] \setminus \{s\}} \subseteq \langle u_1, \dots, u_{b'} \rangle$ for some $s \in [b]$, where $\gamma_s \neq 0$ (obviously, $b-1 \leq b' \leq b$). Each node u_ℓ computes

$$\sum_{t=1}^b \gamma_t \varphi \begin{pmatrix} u_\ell \\ e_t \end{pmatrix} \cdot x = \varphi \begin{pmatrix} u_\ell \\ \sum_{t=1}^b \gamma_t e_t \end{pmatrix} \cdot x = \varphi \begin{pmatrix} u_\ell \\ v_j \end{pmatrix} \cdot x, \quad (1)$$

and sends it to the newcomer.

Notice that the elements

$$\left\{ \varphi \begin{pmatrix} u_\ell \\ e_t \end{pmatrix} \cdot x \right\}_{t \in [b] \setminus r(u_\ell)}$$

were sent to u_ℓ by s in the initial stage (Subsection 3.1). If needed, $\varphi \begin{pmatrix} u_\ell \\ e_{r(u_\ell)} \end{pmatrix} \cdot x$ may be computed using Lemma 2. Hence, every node u_ℓ is capable of performing the computation in (1).

²e.g., a sparse representation of $x = (\gamma_1, \dots, \gamma_B)$ is $\{(j, \gamma_j)\}_{j|\gamma_j \neq 0}$. This representation clearly requires $O(w_H(x) \cdot \log B) = O(w_H(x) \cdot \log b)$ space, where $w_H(x)$ is the Hamming weight of x .

Lemma 4. *By using the information received from Algorithm 1, the newcomer may restore the information from the failed node v_j by using $O(B^2)$ field operations.*

Proof. The newcomer may retrieve $M_{v_j} \cdot x$, the lost information of v_j , by using Lemma 1, Lemma 2, and Lemma 3. Since $\langle e_j \rangle_{j \in [b] \setminus \{s\}} \subseteq \langle u_1, \dots, u_{b'} \rangle$, it follows that the matrix

$$\begin{pmatrix} \varphi(u_1) \\ \vdots \\ \varphi(u_{b'}) \end{pmatrix}$$

has a submatrix which is equivalent to

$$A \triangleq - \begin{pmatrix} \varphi(v_j) \\ \vdots \\ \varphi(v_j) \end{pmatrix},$$

where $\{i_1, \dots, i_{b'}\} = [b] \setminus \{s\}$. By Lemma 3, $\langle A \rangle = P_{\langle v_j \rangle}$ (see Definition 2), and hence A is row equivalent to $-M_{v_j}$. Therefore, the vector $M_{v_j} \cdot x$ may be extracted from the received information by using Observation 1.

Assuming the identity of v_j is known, this algorithm requires communicating either $b-1$ or b field elements. For functional repair no further computations are required. For exact repair the newcomer needs to perform Gaussian-like process on a matrix of size $b' \times B$. By Lemma 1, this process requires the same $O(b^2)$ row operations preformed during a Gaussian elimination of a $b' \times b$ matrix. However, these row operations are being preformed on rows of length B , and hence this Gaussian elimination requires $O(b^2 \cdot B) = O(B^2)$ field operations. \square

Notice that the only requirement on the nodes $u_1, \dots, u_{b'}$ participating in v_j 's repair process is that $\langle e_j \rangle_{j \in [b] \setminus \{s\}} \subseteq \langle u_1, \dots, u_{b'} \rangle$. It follows that if $u_1, \dots, u_{b'}$ are active nodes that form a basis to \mathbb{F}_q^b (i.e. $b' = b$), then it is possible to repair *any* node v_j by using Algorithm 1.

Corollary 1. *Using Algorithm 1, it is possible to add a new node that was not initially in the DSS (see Section 3.1).*

3.3 Local Repair

It is often required that a failed node will be repairable from as few other active nodes as possible. It is clear that without replication of nodes, a minimum of two active nodes is necessary for such a repair. Clearly, such a repair can be done by contacting k nodes from which the reconstruction is possible. In the following we present an alternative repairing approach that may achieve this minimum. The possibility of achieving this minimum depends on the specific assignment of vectors to the nodes. This assignment will be discussed in detail in Section 4.

Algorithm 2. *Let v_j be the failed node and let $\{u_1, \dots, u_\ell\}$ be a set of active linearly independent nodes such that $v_j \in \langle u_1, \dots, u_\ell \rangle$. For all $t \in [\ell]$, the newcomer ν downloads the entire vector $M_{u_t} \cdot x$ from u_t .*

Lemma 5. *By using the information received from Algorithm 2, the newcomer ν may restore the information of the failed node v_j in $O(\ell^2 \cdot b)$ field operations.*

Proof. Since $v_j \in \langle u_1, \dots, u_\ell \rangle$, it follows that $v_j = \sum_{t=1}^\ell \gamma_t u_t$ for some $\gamma_1, \dots, \gamma_\ell \in \mathbb{F}_q$. By the definition of the matrices $\{M_{u_1}, \dots, M_{u_\ell}\}$, ν downloads the set of elements

$$\left\{ \varphi \begin{pmatrix} u_t \\ e_i \end{pmatrix} \cdot x \right\}_{i \in [b] \setminus \{r(v_{i_t})\}}$$

for all $t \in [\ell]$. The missing elements

$$\left\{ \varphi \begin{pmatrix} u_t \\ e_{r(u_t)} \end{pmatrix} \cdot x \right\}_{t \in [\ell]}$$

are computed by Lemma 2 in $O(\ell \cdot b)$ field operations. The newcomer computes the coefficients $\gamma_1, \dots, \gamma_\ell$, e.g. by performing Gaussian elimination on the matrix

$$\begin{pmatrix} u_1 \\ \vdots \\ u_\ell \\ v_j \end{pmatrix},$$

a process requiring $O(\ell^2 \cdot b)$ field operation. Having these coefficients the newcomer performs

$$\sum_{t=1}^\ell \gamma_t \varphi \begin{pmatrix} u_t \\ e_i \end{pmatrix} \cdot x = \varphi \begin{pmatrix} \sum_{t=1}^\ell \gamma_t u_t \\ e_i \end{pmatrix} \cdot x = \varphi \begin{pmatrix} v_j \\ e_i \end{pmatrix} \cdot x.$$

for all $i \in [b] \setminus \{r(v_j)\}$ in $O(\ell \cdot b)$ operations, and reassembles the vector $M_j \cdot x$. Overall, Algorithm 2 requires $O(\ell^2 \cdot b)$ field operations and $\ell \cdot (b - 1)$ communication units. \square

Corollary 2. *Let v_j be a failed node. If ℓ is the smallest integer such v_j is in the linear span of ℓ other active nodes, then the locality of repairing v_j is ℓ .*

3.4 Parallel Repair

Consider the scenario of multiple simultaneous node failures. Obviously, under t failures, if the conditions of Algorithm 2 are satisfied, then it is possible to execute t sequential instances of the repair algorithm. We show that this could be improved in a certain special case. This is a simple consequence of Lemma 5.

Lemma 6. *If $\{v_{i_1}, \dots, v_{i_t}\}$ is a set of failed nodes and $\{v_{j_1}, \dots, v_{j_s}\}$ is a set of active linearly independent nodes, of the remaining nodes, such that*

$$\{v_{i_1}, \dots, v_{i_t}\} \subseteq \langle v_{j_1}, \dots, v_{j_s} \rangle,$$

then it is possible to repair all failures by communicating $s \cdot (b - 1)$ field elements.

Proof. Assume that a third party Ψ is managing the repair process of all t nodes simultaneously. Ψ may download the entire content of all nodes $\{v_{j_1}, \dots, v_{j_s}\}$, and compute the set $\{\varphi \begin{pmatrix} v_{i_m} \\ e_\ell \end{pmatrix} \cdot x\}_{\ell=1}^b$ for each $m \in [t]$ using Algorithm 2. \square

The complexity of Lemma 6 remains t times the complexity of Algorithm 2. However, the amount of communication is the same as in a single instance of Algorithm 2. It is evident that this algorithm requires good locality. An assignment of vectors to nodes that achieves locality is discussed in Subsection 4.2.

3.5 Reconstruction

This subsection presents two reconstruction algorithms for two different models of communication. In Algorithm 3, which follows, the DC accesses b active nodes and downloads their data in its entirety for the reconstruction. The number of communicated field elements is $b(b-1) = 2B$. Algorithm 4, which follows, uses the additional assumption that the nodes participating in the reconstruction know the identities of one another (e.g., by broadcast, shared memory or by acknowledgement from the DC), and guarantees reconstruction by communicating B field elements. This is the minimum communication that guarantees a complete reconstruction of x .

Algorithm 3. Let $\{u_1, \dots, u_b\}$ be a set of active linearly independent nodes. For each $j \in [b]$, the DC downloads the vector $M_{u_j} \cdot x$ from u_j and computes the missing element $\varphi_{(e_r(u_j))}^{(u_j)} \cdot x$ from each node by using Lemma 2. The DC assembles the vector $w \in \mathbb{F}_q^{b^2}$ such that³ $w_{(i,j)} = \varphi_{(e_i)}^{(u_j)} \cdot x$, and the $b^2 \times B$ matrix A whose rows are $\{\varphi_{(e_j)}^{(u_i)}\}_{i,j \in [b]}$. The vector x is then reconstructed by solving the linear system of equations $Ax = w$.

Lemma 7. The matrix A in Algorithm 3 has full rank. In particular, the DC may extract x using $O(B^3)$ field operations.

Proof. By Lemmas 1 and 2, for each $t \in [b]$ the submatrix

$$\begin{pmatrix} \varphi_{(e_t)}^{(u_1)} \\ \vdots \\ \varphi_{(e_t)}^{(u_b)} \end{pmatrix}$$

is row equivalent to the matrix

$$\begin{pmatrix} \varphi_{(e_t)}^{(e_1)} \\ \vdots \\ \varphi_{(e_t)}^{(e_b)} \end{pmatrix}.$$

The matrix A is therefore row equivalent (up to redundant rows) to a matrix whose rows are $\{\varphi_{(e_j)}^{(e_i)}\}_{i \neq j}$, which may clearly be seen as equivalent to the identity matrix of size $B \times B$. Thus, the DC may use Observation 1 to recover x . Computing the rows of A requires $O(b^2B) = O(B^2)$ operations. Solving a $b^2 \times B$ linear system of equations requires additional $O(B^3)$ operations. \square

Assuming that every node participating in the reconstruction algorithm knows the identity of all other participating nodes, it is possible to reduce the communication to merely $|x| = B$ field elements from $b-1$ nodes, $\frac{b}{2}$ elements from each node if b is even and either $\frac{b-1}{2}$ elements or $\frac{b+1}{2}$ elements if b is odd. As mentioned earlier, this is the minimum possible communication since no outer code is used. The following matrix, whose construction is deferred to Appendix A, will be used in Algorithm 4.

Definition 3. Let N be a $b \times b$ matrix over \mathbb{F}_2 such that

- (1) For all $i \in [b]$, $N_{i,b} = 0$.
- (2) For all $i \in [b-1]$, $N_{b,i} = 1$.
- (3) For all $i \in [b]$, $N_{i,i} = 0$.

³The entries of the vector $w \in \mathbb{F}_q^{b^2}$ are identified by the elements of $[b]^2$ according to the lexicographic order.

(4) For all $i, j \in [b], i \neq j$, $N_{i,j} \neq N_{j,i}$.

(5) If b is even then for all $i \in [b-1]$ the Hamming weight of the i th column is $\frac{b}{2}$.

(6) If b is odd then for all $i \in [b-1]$ the Hamming weight of the i th column is either $\frac{b-1}{2}$ or $\frac{b+1}{2}$ and the total Hamming weight of N is $\binom{b}{2}$.

Example 1. The following matrices satisfy the requirements of Definition 3 for $b = 6$ and $b = 5$:

$$\begin{pmatrix} 0 & 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 1 & 0 \end{pmatrix}, \begin{pmatrix} 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 1 & 1 & 1 & 1 & 0 \end{pmatrix}$$

Algorithm 4. Let $\{u_1, \dots, u_b\}$ be any set of linearly independent nodes and let N be the matrix from definition 3. For all $i \in [b-1]$, the DC downloads from node u_i all elements $\varphi_{u_j}^{(u_i)} \cdot x$ such that $N_{j,i} = 1$. The DC assembles the vector $w \in \mathbb{F}_q^{\binom{b}{2}}$ such that⁴ $w_{\{i,j\}} = \varphi_{u_j}^{(u_i)} \cdot x$, and a $B \times B$ matrix A whose rows are the vectors $\{\varphi_{u_j}^{(u_i)}\}_{i \neq j}$. The vector x is then reconstructed by solving the linear system of equations $Ax = w$.

Lemma 8. Using the information received from Algorithm 4, the DC may construct the vector w . In addition, the matrix A from Algorithm 4 has full rank, and hence x is reconstructible by using $O(B^3)$ field operations.

Proof. By (4) of Definition 3 it is evident that for all $i, j \in [b], i \neq j$, the DC receives the element $\varphi_{u_j}^{(u_i)} \cdot x$ exactly once. To prove that the reconstruction of x is possible, we show that A is row equivalent to a matrix whose rows are $\{\varphi_{e_j}^{(e_i)}\}_{i \neq j}$. The latter may easily be seen as equivalent to the $B \times B$ identity matrix. For any $i \in [b]$, add the zero row $\varphi_{u_i}^{(u_i)}$ to the proper submatrix to get

$$\begin{pmatrix} \varphi_{u_1}^{(u_i)} \\ \vdots \\ \varphi_{u_b}^{(u_i)} \end{pmatrix}.$$

By Lemma 1 this matrix is row equivalent to

$$\begin{pmatrix} \varphi_{e_1}^{(u_i)} \\ \vdots \\ \varphi_{e_b}^{(u_i)} \end{pmatrix}.$$

By rearranging the rows of A we may consider submatrices of the form

$$\begin{pmatrix} \varphi_{e_i}^{(u_1)} \\ \vdots \\ \varphi_{e_i}^{(u_b)} \end{pmatrix}.$$

⁴The entries of the vector $w \in \mathbb{F}_q^{\binom{b}{2}}$ are identified by all 2-subsets of $[b]$ according to the lexicographic order.

for all $i \in [b]$. These submatrices are row equivalent by Lemma 1 to

$$\begin{pmatrix} \varphi_{e_i}^{(e_1)} \\ \vdots \\ \varphi_{e_i}^{(e_b)} \end{pmatrix}.$$

Hence, by addition of redundant rows, we get that A is equivalent to the identity matrix. Thus, the reconstruction of x is possible by Gaussian elimination, requiring $O(B^3)$ operations. \square

3.6 Modification

A useful property of a DSS is being able to update a small fraction of x without having to initialize the entire system. The linear nature of our code and the absence of an outer code allows these modifications to be done efficiently. In particular, the complexity of the process is a function of the Hamming distance $d_H(x, y)$, where y is the modification of the vector x . In MDS based distributed storage systems a change of a single bit of x usually requires changing a large portion of the data. Therefore, one more advantage of our system is revealed.

Lemma 9. *If $x \in \mathbb{F}_q^B$ is stored in the system, it is possible to update the system to contain $y \in \mathbb{F}_q^B$ by communicating $(\log B + \log q) \cdot d_H(x, y) \cdot n$ bits.*

Proof. Each node receives a list $\{(\delta_i, \ell_i)\}_{i=1}^{d_H(x, y)}$, where $\delta_i \in \mathbb{F}_q$ and $\ell_i \in [B]$. The list indicates the values of the nonzero entries of the vector $y - x$. Each node v , holding the vector $M_v \cdot x$ (see Section 3.1) may assemble the matrix M_v and compute:

$$M_v \cdot x + M_v \cdot (y - x) = M_v \cdot y.$$

Communicating the list $\{(\delta_i, \ell_i)\}_{i=1}^{d_H(x, y)}$ to all the n nodes clearly requires $(\log B + \log q) \cdot d_H(x, y) \cdot n$ bits. \square

4 Assignment of Vectors

In Section 3 we proved that the performance of the detailed algorithms strongly relies on the chosen vectors v_1, \dots, v_n . Since both repair and reconstruction algorithms require linearly independent nodes, it follows that the assigned set of vectors should contain a basis to \mathbb{F}_q^b even after multiple failures.

Choosing $n = \begin{bmatrix} b \\ 1 \end{bmatrix}_q$ and assigning *all* possible normalized vectors would suffice for repairing exponentially many failures. However, using $\begin{bmatrix} b \\ 1 \end{bmatrix}_q = \Theta(q^b)$ storage nodes to store a file of size $B = \Theta(b^2)$ is unnecessary, as will be shown in the sequel. Furthermore, expecting exponentially many failures is nonrealistic.

In order to achieve reasonable failure resilience using a reasonable number of nodes, it suffices to consider the case $n = O(b)$. Subsection 4.1 discusses an assignment of vectors compatible with Algorithm 1 presented in Subsection 3.2. An assignment compatible with Algorithms 2 of Subsection 3.3 and also for the algorithm of Subsection 3.4 is presented in Subsection 4.2.

Definition 4. *For $t \in \mathbb{N}$ a set $S \subseteq \mathbb{F}_q^b$ is called a t -resilient spanning set if every t -subset T of S satisfies $\langle S \setminus T \rangle = \mathbb{F}_q^b$.*

Observation 2. If S is a t -resilient spanning set then by using $|S|$ storage nodes assigned with the vectors in S (see Subsection 3.1) then it is possible to repair and reconstruct in the presence of up to t simultaneous node failures.

Example 2. The following set is a 2-resilient spanning set in \mathbb{F}_2^7 :

$$\begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 0 & 1 & 0 & 1 \end{pmatrix}$$

4.1 Minimum Bandwidth Assignment

In what follows we present a construction of a set of vectors $\{v_1, \dots, v_n\}$ compatible with Algorithm 1, achieving $d\beta \leq \alpha + 1$.

Lemma 10. Let $b \in \mathbb{N}$ and let C be a linear block code of length $c \cdot b$ for some constant $c > 0$, dimension b , and minimum Hamming distance δ over \mathbb{F}_q . If M is a generator matrix of C then the columns of M are a $(\delta - 1)$ -resilient spanning set (see Definition 4).

The inverse of Lemma 10 is also true, as stated in the next lemma.

Lemma 11. Let $S \subseteq \mathbb{F}_q^b$ be an assignment of vectors to nodes in some DSS which is resilient to t node failures by using the algorithms described in Section 3. If G is the matrix whose columns are the elements of S and $C \triangleq \{xG \mid x \in \mathbb{F}_q^b\}$ then C is a linear code of minimum Hamming distance $t + 1$.

Example 3. Let C to be a binary Justesen code [17] of length $O(b)$, dimension b , and minimum Hamming distance δb . We get that the corresponding code (see Section 3) uses $O(b)$ storage nodes while being able to recover from any δb simultaneous node failures. In addition, the code uses the binary field. This choice admits the following parameters: $q = 2$, $B = \binom{b}{2}$, $n = O(B^{1/2})$, $d = b = O(B^{1/2})$, $k = b = O(B^{1/2})$, $\alpha = b - 1 = O(B^{1/2})$, and $\beta = 1$.

4.2 Minimum Locality Assignment

Algorithm 2 in Subsection 3.3 may possibly achieve the optimal locality. It is evident from Lemma 5 that in order to get good locality, the set $\{u_1, \dots, u_\ell\}$ from Algorithm 2 is required to be small. However, this requirement conflicts with the requirements of Algorithms 1, 3, and 4, since they all involve large linearly independent sets.

In this subsection we show that by choosing some basis of \mathbb{F}_q^b , partitioning it to equally sized subsets and taking the linear span of each subset, some locality is achievable. The resulting failure resilience will grow with the field size. Thus, this technique will be particularly useful in large fields.

Definition 5. Let c be a positive integer such that c divides b , and let $A \triangleq \{v_1, \dots, v_b\}$ be a basis of \mathbb{F}_q^b . Partition A into $\frac{b}{c}$ equally sized subsets $A_i \triangleq \{v_{ic+1}, \dots, v_{(i+1)c}\}$ for $i \in \{0, \dots, \frac{b}{c} - 1\}$. Let $V_i \subseteq \mathbb{F}_q^b$ be a set of $\binom{c}{1}_q$ representatives for the 1-subspaces of $\langle A_i \rangle$. Finally, let $V \triangleq \bigcup_{i=1}^{\frac{b}{c}} V_i$.

Lemma 12. *The set V from Definition 5 is a $(q^{c-1} - 1)$ -resilient spanning set (see Definition 4). Furthermore, assigning V to nodes in a DSS allows repairing any node failure using at most c active nodes.*

Proof. Since

$$q^{c-1} - 1 = \begin{bmatrix} c \\ 1 \end{bmatrix}_q - \begin{bmatrix} c-1 \\ 1 \end{bmatrix}_q - 1 < \begin{bmatrix} c \\ 1 \end{bmatrix}_q - \begin{bmatrix} c-1 \\ 1 \end{bmatrix}_q,$$

it follows that after any set of at most $q^{c-1} - 1$ node failures, the set of remaining active nodes in any V_i is not contained in any $(c-1)$ -subspace of $\langle A_i \rangle$. Therefore, any V_i still contains a basis for $\langle A_i \rangle$. Since $\langle A_1 \rangle \oplus \dots \oplus \langle A_{b/c} \rangle = \mathbb{F}_q^b$, it follows that V is $(q^{c-1} - 1)$ -resilient spanning set.

Let v_j be a failed node and let V_t be the set containing it. We have to prove that v_j is repairable using at most c other nodes in the presence of at most $q^{c-1} - 1$ failures. We have shown that after $q^{c-1} - 1$ failures, the remaining active nodes in any given V_i contain a basis of $\langle A_i \rangle$. Let $\{u_1, \dots, u_c\} \subseteq \langle A_t \rangle$ be such a basis in V_t . It follows that $v_j \in \langle u_1, \dots, u_c \rangle$, and hence v_j is repairable by accessing at most c nodes by Lemma 5. \square

This construction requires $\frac{b}{c} \cdot \begin{bmatrix} c \\ 1 \end{bmatrix}_q$ nodes and allows locality of c in the presence of up to $q^{c-1} - 1$ failures. For simple comparison, the trivial replication code with $\frac{b}{c} \cdot \begin{bmatrix} c \\ 1 \end{bmatrix}_q$ nodes allows locality of 1 in the presence of up to $\frac{1}{c} \cdot \begin{bmatrix} c \\ 1 \end{bmatrix}_q - 1$ failures. We note that

$$\frac{q^{c-1} - 1}{\frac{1}{c} \cdot \begin{bmatrix} c \\ 1 \end{bmatrix}_q - 1} \xrightarrow{q \rightarrow \infty} c,$$

and in particular for $c = 2$,

$$\frac{q^{2-1} - 1}{\frac{1}{2} \cdot \begin{bmatrix} 2 \\ 1 \end{bmatrix}_q - 1} = 2.$$

Therefore, this code outperforms the trivial one by approximately a factor of c for large field size, while providing low locality. In particular, a minimal locality of 2 is achievable for any q .

References

- [1] C. Suh and K. Ramchandran, “Exact-repair MDS code construction using interference alignment,” *IEEE Transactions on Information Theory*, vol. 57, no. 3, pp. 1425–1442, 2011.
- [2] F. Oggier and A. Datta, “Self-repairing homomorphic codes for distributed storage systems,” in *Proceedings of INFOCOM*, pp. 1215–1223, 2011.
- [3] Z. Wang, I. Tamo, and J. Bruck, “Long MDS codes for optimal repair bandwidth,” in *Proceedings IEEE International Symposium on Information Theory (ISIT)*, pp. 1182–1186, 2012.
- [4] I. Tamo, Z. Wang, and J. Bruck, “Access vs. bandwidth in coded for distributed storage,” *arxiv:1303.3668*, 2014.
- [5] D. S. Papailiopoulos, A. G. Dimakis, and V. R. Cadambe, “Repair optimal erasure codes through Hadamard designs,” in *49th Annual Allerton Conference on Communication, Control, and Computing*, pp. 1382–1389, 2011.

- [6] N. B. Shah, K. Rashmi, P. V. Kumar, and K. Ramchandran, “Interference alignment in regenerating codes for distributed storage: Necessity and code constructions,” *IEEE Transactions on Information Theory*, vol. 58, no. 4, pp. 2134–2158, 2012.
- [7] N. Silberstein, A. S. Rawat, and S. Vishwanath, “Error resilience in distributed storage via rank-metric codes,” in *50th Annual Allerton Conference on Communication, Control, and Computing*, pp. 1150–1157, 2012.
- [8] N. Silberstein, A. Rawat, O. Koyluoglu, and S. Vishwanath, “Optimal locally repairable codes via rank-metric codes,” in *Proceedings IEEE International Symposium on Information Theory (ISIT)*, pp. 1819–1823, 2013.
- [9] K. V. Rashmi, N. B. Shah, and P. V. Kumar, “Optimal exact-regenerating codes for distributed storage at the MSR and MBR points via a product-matrix construction,” *IEEE Transactions on Information Theory*, vol. 57, no. 8, pp. 5227–5239, 2011.
- [10] A. Dimakis, P. Godfrey, Y. Wu, M. Wainwright, and K. Ramchandran, “Network coding for distributed storage systems,” *IEEE Transactions on Information Theory*, vol. 56, no. 9, pp. 4539–4551, 2010.
- [11] F. Oggier and A. Datta, “Self-repairing codes for distributed storage - A projective geometric construction,” in *Information Theory Workshop (ITW)*, pp. 30–34, 2011.
- [12] H. Hollmann, “Storage codes; coding rate and repair locality,” in *International Conference on Computing, Networking and Communications (ICNC)*, pp. 830–834, 2013.
- [13] T. Etzion and N. Raviv, “Equidistant codes in the Grassmannian,” *arXiv:1308.6231v2 [math.CO]*, 2014.
- [14] J. H. van Lint and R. M. Wilson, *A course in combinatorics*. Cambridge university press, 2001.
- [15] A. Beutelspacher and U. Rosenbaum, *Projective geometry: from foundations to applications*. Cambridge University Press, Cambridge, 1998.
- [16] A. A. Mofrad, M.-R. Sadeghi, and D. Panario, “Solving sparse linear systems of equations over finite fields using bit-flipping algorithm,” *Linear Algebra and its Applications*, vol. 439, no. 7, pp. 1815 – 1824, 2013.
- [17] J. Justesen, “Class of constructive asymptotically good algebraic codes,” *IEEE Transactions on Information Theory*, vol. 18, no. 5, pp. 652–656, 1972.

Appendix A

Two constructions of a matrix satisfying the requirements of Definition 3 are given, Construction 1 for even b and Construction 2 for odd b . It is easily verified that these two constructions satisfy the requirement of Definition 3.

Construction 1. Let b be an even integer. Define $N \in \mathbb{F}_2^{b \times b}$ as follows. For all $i \in [b]$ let $N_{i,b} = 0$ and for all $i \in [b-1]$ let $N_{b,i} = 1$. The remaining $(b-1) \times (b-1)$ upper left submatrix is defined as follows. The first row is the $b-1$ bit vector $0^{b/2}1^{b/2-1}$; that is, $\frac{b}{2}$ zeros followed by $\frac{b}{2}-1$ ones. The rest of the rows are all cyclic shifts of it (see Example 1 for the case $b=6$).

Construction 2. Let b be an odd integer. Define $N \in \mathbb{F}_2^{b \times b}$ as follows. For all $i \in [b]$ let $N_{i,b} = 0$ and for all $i \in [b-1]$ let $N_{b,i} = 1$. The remaining $(b-1) \times (b-1)$ upper left submatrix is defined as follows. The first row is the $b-1$ bit vector $0^{(b+1)/2}1^{(b-3)/2}$, that is, $\frac{b+1}{2}$ zeros followed by $\frac{b-3}{2}$ ones. The rest of the rows are all cyclic shifts of it. In addition, set the sub diagonal entries $(1, \frac{b-1}{2} + 1), (2, \frac{b-1}{2} + 2), \dots, (\frac{b-1}{2}, b-1)$ to 1 (see Example 1 for the case $b = 5$).